

Date: 27/07/2024

WEEK 2

## Optimization Algorithms

⇒ Batch v/s Mini-batch gradient descent

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(m)}]$$

(n, m)

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

(l, m)

What if  $m = 5\text{million}$  (5,000,000)

mini-batches of 1000 each

$$X = \left[ \underbrace{x^{(1)} \rightarrow x^{(1000)}}_{\{n, 1000\}} \mid x^{(1001)} \rightarrow x^{(2000)} \dots \right] \quad \left[ \underbrace{x^{(1001)} \rightarrow x^{(2000)} \dots}_{\{n, 1000\}} \right]$$

$\{5000\}$

Similarly for  $Y$  too;  
(1, 1000)

Date:

Mini-batch  $t$ :  $x^{\{t\}}, y^{\{t\}}$

⇒ Mini-batch gradient descent :

Vectorised:

for  $t = 1, \dots, 5000$

{

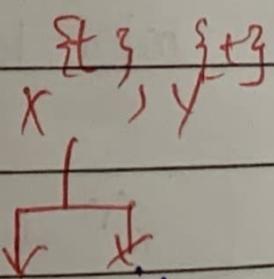
Forward prop on  $x^{\{t\}}$

$$z^{[1]} = w^{[1]} x^{\{t\}} + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

}

$$a^{[L]} = g^{[L]}(z^{[L]})$$



Compute cost  $J = \frac{1}{1000} \sum_{i=1}^l \mathcal{L}(y^{(aci)}, \hat{y}^{(i)})$

$$+ \frac{\lambda}{2 \cdot 1000} \sum_l \|w^{[L]}\|_F^2$$

Regularisation

Date:

Backprop to compute gradients  
wrt  $J^{\{t\}}$  (using  $x^{\{t\}}$ ,  
 $y^{\{t\}}$ )

$$w^{[l]} := w^{[l]} - \alpha d w^{[l]}$$

$$b^{[l]} := b^{[l]} - \alpha d b^{[l]}$$

}

"Epoch"  $\rightarrow$  pass through the  
training set

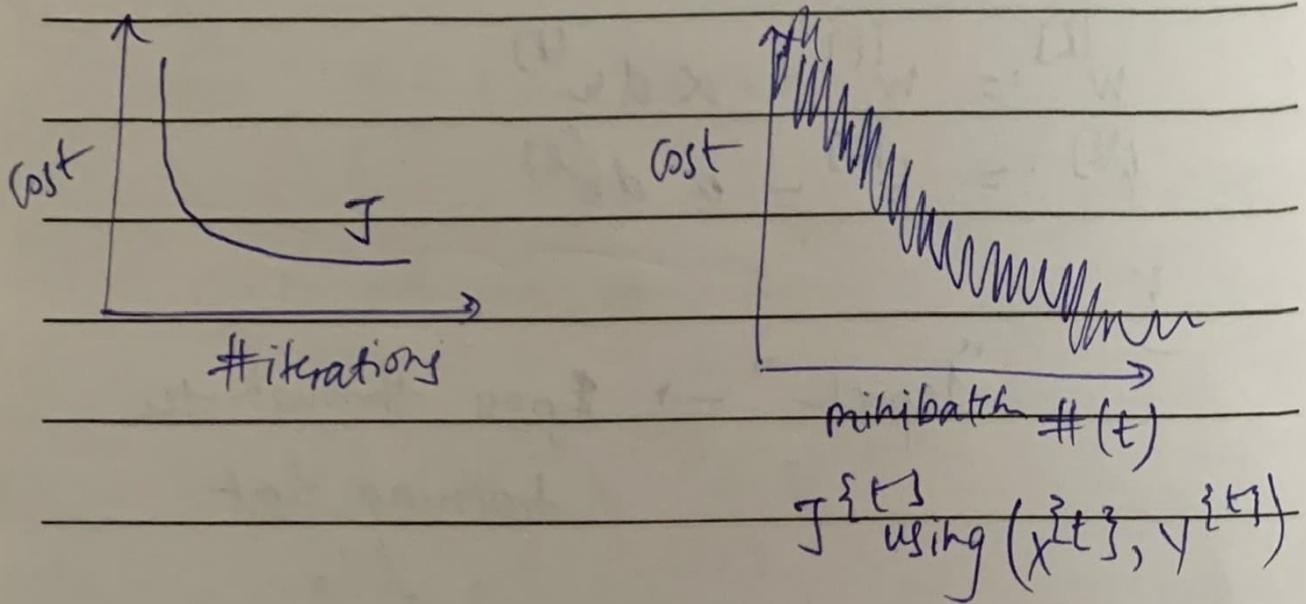
When using a large dataset; the mini-batch  
is works much faster.

Date:

→ Training with mini-batch gradient descent :

Batch gradient descent

Mini-batch grad. descent



Not strictly decreasing,  
but should trend downwards

→ Choosing mini-batch size:

mini-batch size = m ;

Batch grad. descent

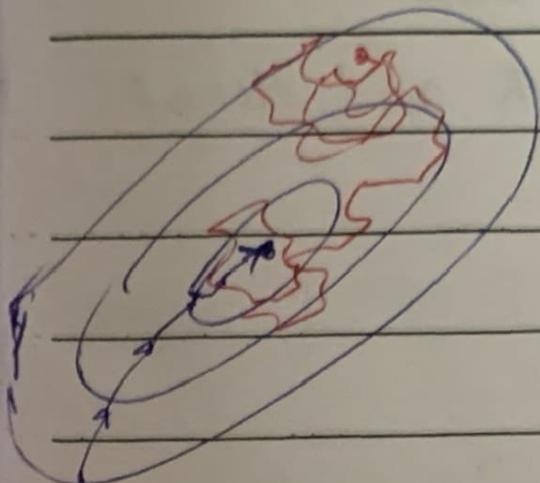


$x^{\{i\}}, y^{\{i\}} = X, Y$

size = 1 ; Stochastic grad. descent

Every example is its own  
mini-batch.

Red for Stochastic (doesn't stop  
Blue for batch at minimum  
Date: always around  
minimum)



In practice:

size in batch

l and m.

① Batch grad. descent.

Too long per iteration

② Stochastic grad. descent

loosing all speed from vectorisation

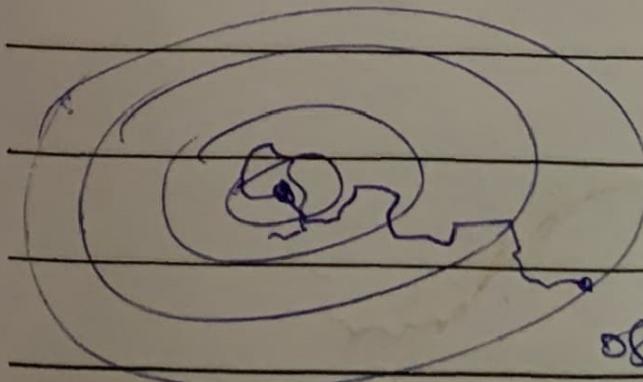
∴ each example is processed: inefficient

③ In batch l and m:

too big/small

→ fastest learning → vectorisation ✓

→ make progress without processing  
the entire training set



not exactly at minimum

(converge) but can be  
oscillated very close to minimum  
by reducing learning rate

Date:

① small training set: Batch grad. descent  
 $(m \leq 2000)$

② typical mini batch sizes:

64, 128, 256, 512 ...

$2^6, 2^7, 2^8, 2^9$  (power of 2)

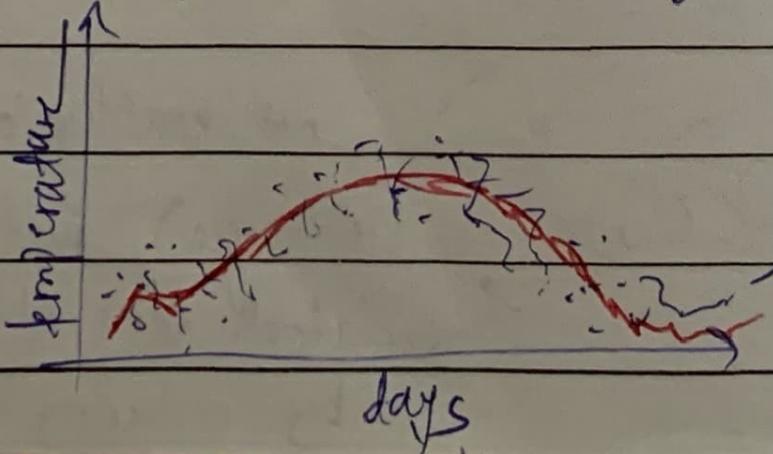
common  $\therefore$  how the computer memory is organised

Make sure mini-batch size fits

~~its~~ in CPU / GPU memory

⇒ Different optimization Algorithms:

① Exponentially weighted moving averages



Date: \_\_\_\_\_

$$V_t = \beta \cdot V_{t-1} + (1-\beta) \theta_t$$

$\beta = 0.9$  :  $\approx 10$  day  
temp.

$V_t$  = approximately  
averaging over

$\beta = 0.98$  :  $\approx 50$  days  
temp.

$\approx \frac{1}{1-\beta}$  days'  
temperature

as  $\beta \uparrow$  days  $\uparrow$  : .  
smooth curves.

$\beta = 0.5$  :  $\approx 2$  days

very noisy, very quickly adapting to  
temperatures

$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{99} = 0.9 V_{98} + 0.1 \theta_{99}$$

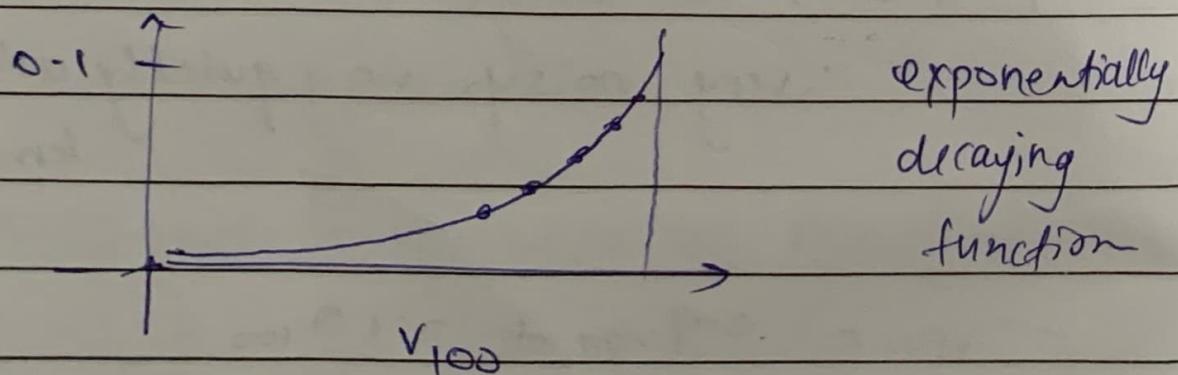
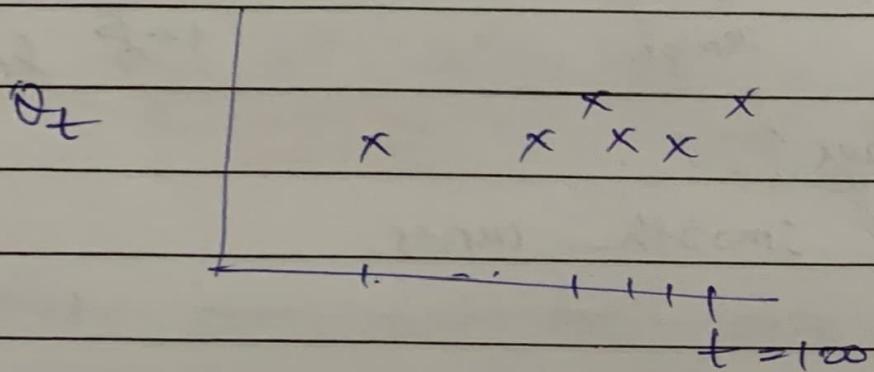
{

$$(0.1 \theta_{99} + 0.9 V_{98})$$

$$\rightarrow V_{100} = 0.1 \theta_{100} + 0.9 \underbrace{V_{99}}_{(0.1 \theta_{99} + 0.9 V_{98})} \dots \text{-- continues.}$$

Date:

$$v_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} \\ + 0.1 \times (0.9)^2 \theta_{98} \\ + 0.1 \times (0.9)^3 \theta_{97} \\ + \dots$$



$$0.9^{10} \approx 0.35 \approx \frac{1}{e}$$

Date:

$$\frac{(1-\varepsilon)}{0.9}^{\frac{1}{\varepsilon}} = \frac{1}{e}$$

means focuses  
on only today.

$$0.98^{\frac{50}{\varepsilon}} \approx \frac{1}{e} \quad \vdash 50 \text{ days when } \beta = 0.98$$

Implementation:

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1-\beta) \alpha_1$$

$$v_2 = \beta v_1 + (1-\beta) \alpha_2$$

$$v_0 = 0$$

$$v_1 = \beta v + (1-\beta) \alpha_1$$

$$v_2 = \beta v + (1-\beta) \alpha_2$$

Date:

$$V_0 = 0$$

Repeat }

Get next  $\theta_t$

$$V_t := \beta V_{t-1} + (1-\beta) \theta_t$$

}

$\Rightarrow$  Bias correction:

Change the formula

$$\frac{V_t^*}{1 - \beta^t}$$

$$\therefore \beta^t \quad (1 - \beta^t)^{-1}$$

$$\hat{V}_t \approx V_t^* \quad \text{at first}$$

two days also it's a good estimate.

Date:

(2) Momentum / Gradient descent with momentum  
(works faster than batch grad. descent)

Momentum :

$V_{dw} = 0$ ;  $V_{db} = 0$  (initialised to matrix of zeros)  
on iteration  $t$ :

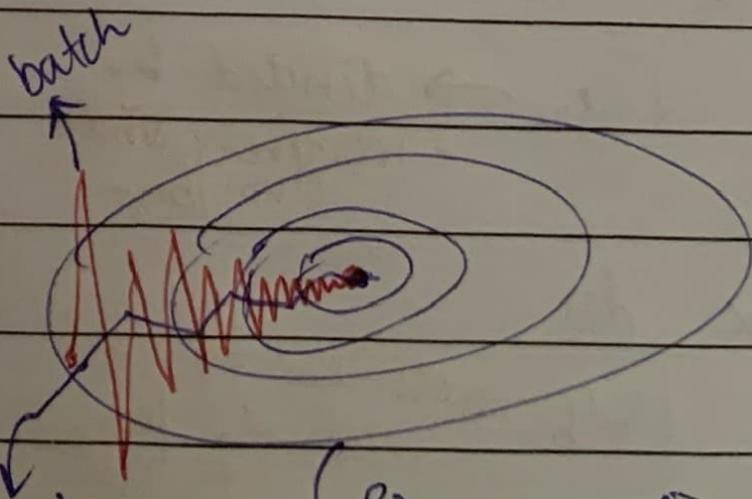
Compute  $dw, db$  on current mini-batch:

$$V_{dw} = \beta V_{dw} + (1-\beta) dw$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$w = w - \alpha V_{dw}$$

$$b = b - \alpha V_{db}$$



Hyperparameters:  $\alpha, \beta$

$$\beta = 0.9$$

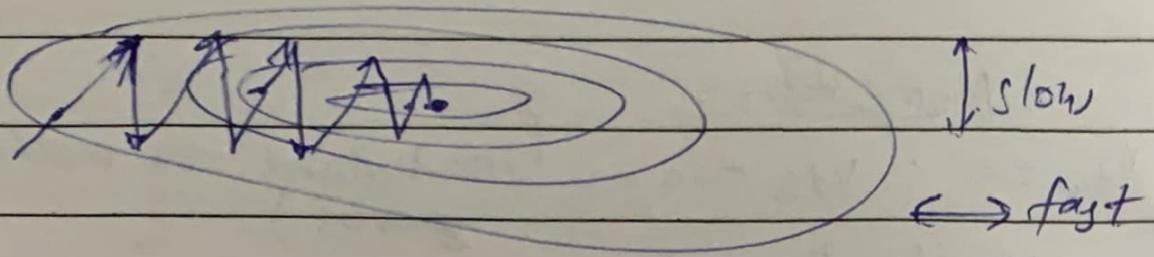
$\approx 10^{previous}$  gradients

average

(Bias correction  
not taken into consideration)

Date:

③ RMS prop  
(root mean squared prop)



On iteration  $t$ :

compute  $dw, db$  on current mini-batch

$$Sdw = \beta_2 Sdw + (1 - \beta_2) dw^2 \quad \text{element-wise}$$

$$Sdb = \beta_2 Sdb + (1 - \beta_2) db^2$$

$$w := w - \alpha \frac{dw}{\sqrt{Sdw}}$$

$\sqrt{Sdw} \rightarrow$  divided by a relatively small number

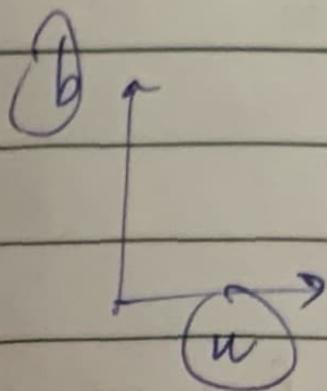
$$b := b - \alpha \frac{db}{\sqrt{Sdb}}$$

$\sqrt{Sdb} \rightarrow$  divided by a relatively large number

Note:  $Sdw + \epsilon$  and

$Sdb + \epsilon$  is used so that  $den^r \neq 0$

Date:



want  $b \downarrow$  and  $w \uparrow$   
in that is happening  
by RMS prop.

A higher learning rate can then be used, for faster learning.

#### ④ Adam optimization algorithm:

Adaptive moment estimation

$$V_{dw}, S_{dw}, V_{db}, S_{db} = 0;$$

On iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw \quad \text{"momentum"}$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2 \quad \text{"RMS prop"}$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \beta_2$$

Date:

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1 t)$$

$$V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1 t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2 t)$$

$$S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2 t)$$

} Bias correction

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw} + \epsilon}}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corra}} + \epsilon}}$$

Commonly used; very effective;

Hyperparameters:

$\alpha$  : needs to be tuned

$\beta_1$  : 0.9 ( $dw$ )

value

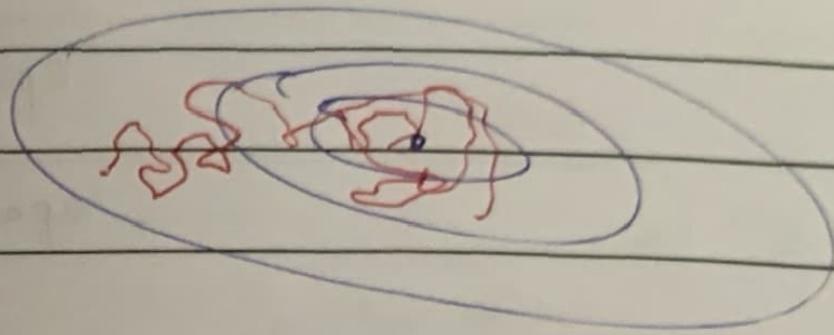
$\epsilon = 10^{-8}$

$\beta_2$  : 0.999 ( $dw^2$ )

value

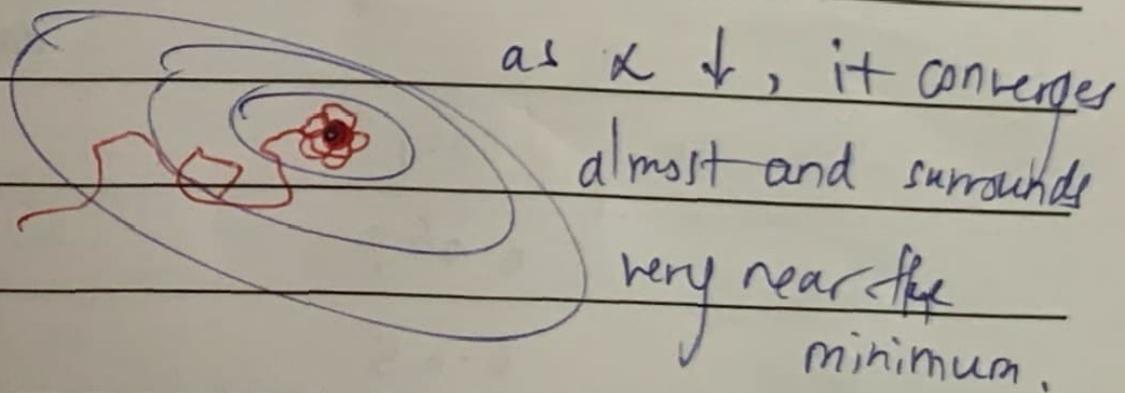
Date:

⇒ Learning rate decay :



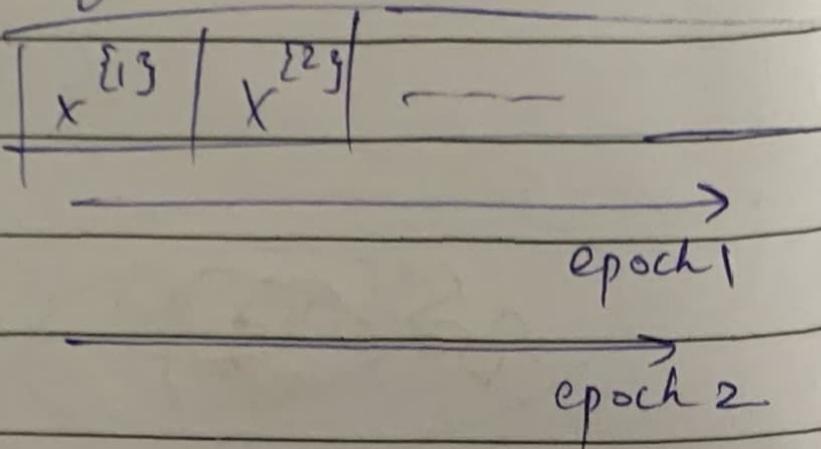
while using mini-batches:

After slowly reducing of  $\alpha$ :



Date: \_\_\_\_\_

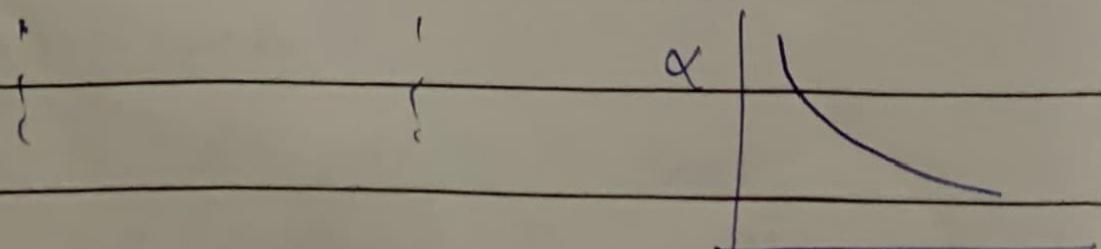
1 epoch =  $\frac{1}{k}$  pass through data



$$\alpha_n = \frac{1}{1 + \text{DecayRate} * \text{epochNumber}}$$

$\alpha_0$

Epoch	$\alpha$	$\alpha_0 = 0.2$
1	0.1	decay-rate = 1
2	0.067	
3	0.05	
4	0.04	



Date:

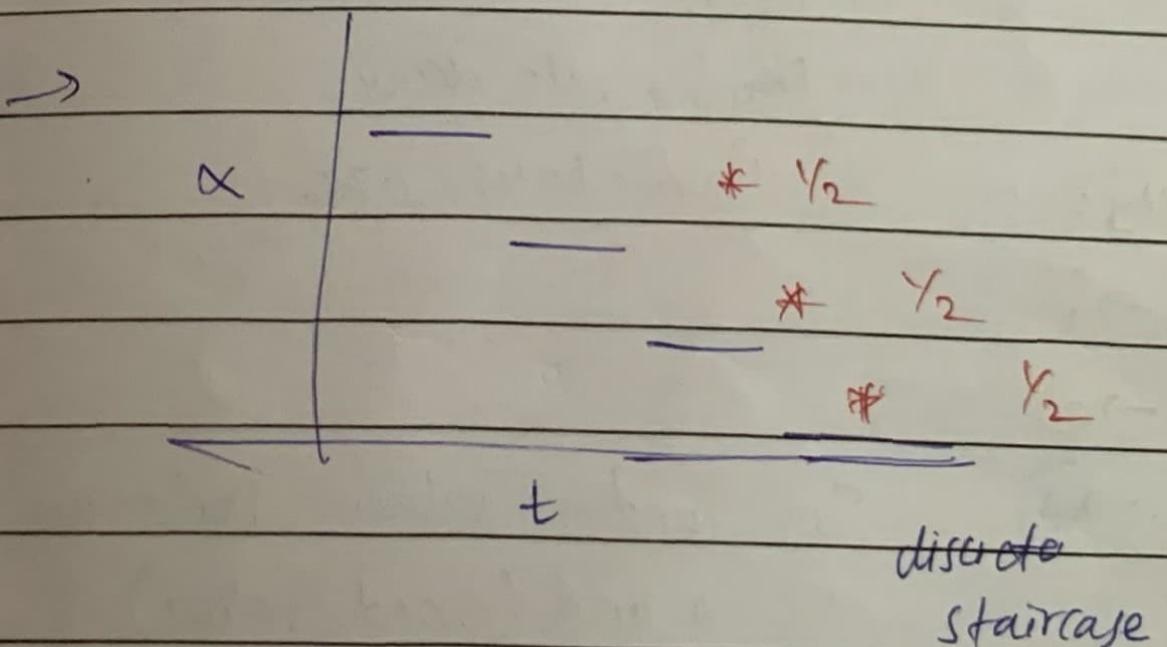
other methods : (for decay rate) :

$$\rightarrow \alpha = 0.95^{\frac{\text{epoch-number}}{\alpha_0}}$$

(exponentially decay)

$$\rightarrow \alpha = \underbrace{\frac{(k)}{\text{epoch-number}}}_{\text{constant}} \cdot \alpha_0 \quad \Rightarrow \quad \frac{k}{\sqrt{E_i}} \cdot \alpha_0$$

$\rightarrow$  mini-batch number



→ Manual decay